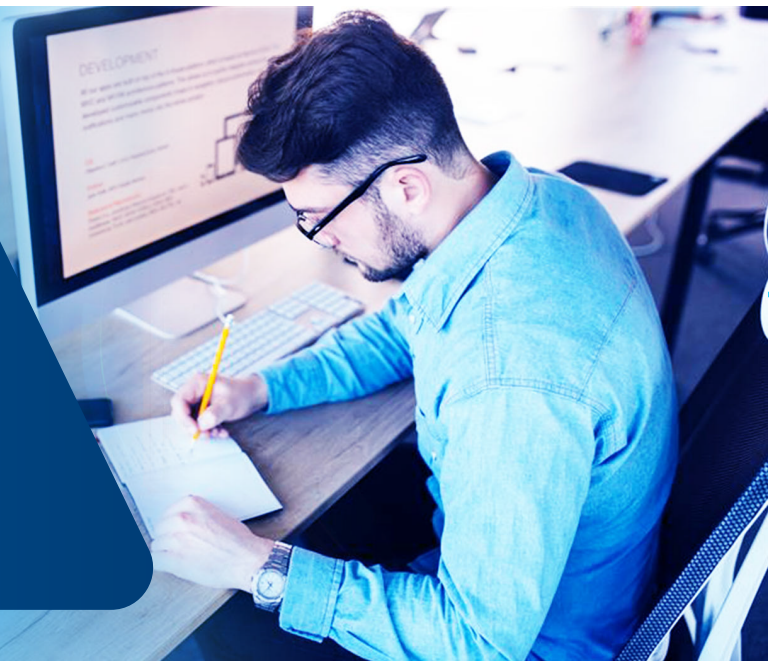


Content Moderation



Overview

Struggling with diverse Indian content? JioCloud Content Moderation uniquely detects abusive or inappropriate material across Hinglish, regional scripts, and mixed formats. Benefit from real-time checks, easy API integration, a dedicated Playground for pre-live testing, and empower your teams with independent, cost-effective, pay-per-review management.

Key Features

- **Multilingual offensive content detection**
Flags explicit, abusive, or hateful content in Indian scripts and mixed-language inputs.
- **Context-aware NLP**
Understands regional slang, sarcasm, and informal phrasing.
- **Real-time flagging**
Detects violations before content is posted or published.
- **Playground for testing and validation**
Test real samples and review labels instantly.
- **API-first architecture**
Add moderation to UGC systems, chat apps, or content platforms easily.
- **Complete self-service**
Manage provisioning, usage, and credits from a unified dashboard.
- **Pay-as-you-go pricing**
Affordable, metered billing with no fixed contract.

Benefits

- Reduce exposure to harmful or abusive content
- Moderate regional and code-mixed text with cultural accuracy
- Lower manual review effort and moderation costs
- Protect brand safety and enable policy compliance
- Improve trust and safety across user communities
- Integrate quickly with APIs and test in a live Playground

Supported Language

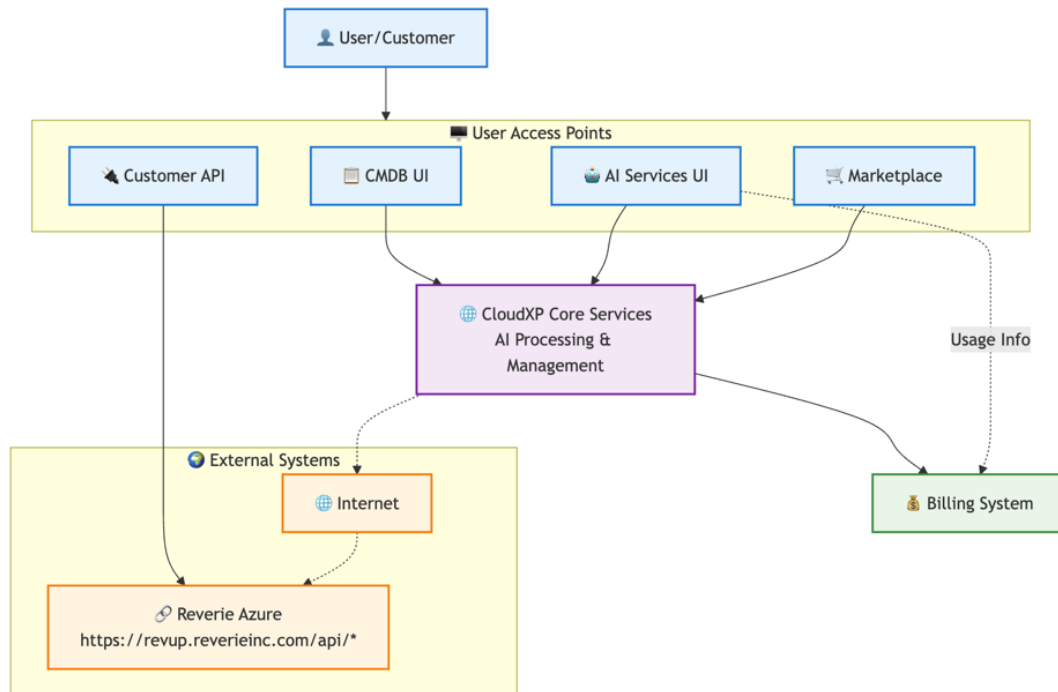
Category	Details
Hindi	Full
Bengali	Full
Tamil	Full
Telugu	Full
Kannada	Full
Malayalam	Full

Category	Details
Marathi	Full
Gujarati	Full
Punjabi	Full
Odia	Full
Assamese	Full
English	Full

Technical Specifications

Category	Details
Moderation Type	Rule-based + ML/NLP-based classification
Categories Supported	Abuse, Hate Speech, Adult, Violence, Custom Terms
Input Format	UTF-8 plain text
Output Format	JSON with category flags and confidence scores
Access Protocol	REST API over HTTPS
Authentication	Bearer Token
Average Latency	2035.4 ms
Rate Limiting	Tiered quotas; adjustable per customer
Web Playground	UI for live testing and validation
Deployment Model	Fully managed SaaS
Billing Model	Per-character consumption-based pricing
Security	HTTPS, zero input/output retention

Architecture Diagram



Use Cases

- Social platforms and messaging apps**
 Flag offensive posts or messages in real time to reduce hate speech and build safer communities.
- eCommerce reviews and Q&A**
 Filter abusive or fake reviews to protect buyer experience and seller reputation.
- Gaming and live chat environments**
 Detect harassment and toxic behavior in fast-moving chats across multiple languages.